

日 本 国 特 許 庁

PATENT OFFICE
JAPANESE GOVERNMENT

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日

Date of Application:

1999年12月28日

出 願 番 号

Application Number:

平成11年特許願第373272号

出 願 人

Applicant(s):

松下電器産業株式会社

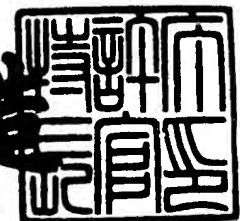


CERTIFIED COPY OF
PRIORITY DOCUMENT

2000年12月22日

特許庁長官
Commissioner,
Patent Office

及 川 耕 造



出証番号 出証特2000-3108192

【書類名】 特許願

【整理番号】 2030714055

【提出日】 平成11年12月28日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 17/27

【発明者】

【住所又は居所】 大阪府門真市大字門真 1 0 0 6 番地
松下電器産業株式会社内

【氏名】 飯塚 泰樹

【特許出願人】

【識別番号】 000005821

【氏名又は名称】 松下電器産業株式会社

【代理人】

【識別番号】 100082692

【弁理士】

【氏名又は名称】 蔵合 正博

【電話番号】 03(3519)2611

【手数料の表示】

【予納台帳番号】 013549

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9004843

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 単語分割方法と装置

【特許請求の範囲】

【請求項 1】 文を単語に分割する単語分割方法において、単語分割されていない文書データから文字結合度を文字接続確率という形で統計的に計算し、計算した文字接続確率を分割対象の文に適用し、接続確率が低い部分で分割することで文を単語に分割することを特徴とする単語分割方法。

【請求項 2】 単語分割されていない文書データから文字結合度を文字接続確率という形で統計的に計算する際、この文字結合度として、複数文字からなるある文字列が出現した後にある文字が出現する確率として計算し、この確率を分割対象の文に適用し、接続確率が低い部分で分割することで文を単語に分割することを特徴とする請求項 1 記載の単語分割方法。

【請求項 3】 単語分割されていない文書データから文字結合度を文字接続確率という形で統計的に計算する際、この文字結合度として、複数文字からなるある文字列が出現した後に複数文字からなるある文字列が出現する確率として計算し、この確率を分割対象の文に適用し、接続確率が低い部分で分割することで文を単語に分割することを特徴とする請求項 1 記載の単語分割方法。

【請求項 4】 文書データから文字結合度としての文字接続確率を統計的計算する際、複数文字からなるある文字列が出現した後にある文字が出現する確率と複数文字からなるある文字列が出現する前にある文字が出現する確率とから、複数文字からなるある文字列が出現した後に複数文字からなるある文字列が出現する確率を計算することを特徴とする請求項 3 記載の単語分割方法。

【請求項 5】 文を単語に分割する単語分割方法において、単語分割されていない学習用文書データから文字結合度としての文字接続確率を統計的に計算し、計算した文字接続確率を分割対象の文に適用し、分割対象の文の中に文字接続確率として計算した以外の文字組み合わせが出現した場合は、前記計算文字接続確率から目的とする確率を推定し、接続確率が低い部分で分割することで文を単語に分割することを特徴とする単語分割方法。

【請求項 6】 文字接続確率による単語分割において、分割するかどうかの判

断の基準となる確率値の閾値を、分割後の平均単語長をもとに動的に決定することを特徴とする請求項 1 から 5 のいずれかに記載の単語分割方法。

【請求項 7】 日本語の文の単語分割において、文字接続確率の他に、漢字や平仮名といった文字の種類が変化する点で単語が分割されやすいという特徴を併用することで単語分割点を決定することを特徴とする請求項 1 から 6 のいずれかに記載の単語分割方法。

【請求項 8】 文字接続確率の他に、カッコなどの記号部分では必ず単語が分割されるということを併用することで単語分割点を決定することを特徴とする請求項 1 から 7 のいずれかに記載の単語分割方法。

【請求項 9】 文を単語に分割する単語分割装置において、文の集りである文書データを入力する文書入力手段と、前記入力された文書データを蓄える文書データ蓄積手段と、前記蓄えられた文書データから文字接続確率を計算する文字接続確率計算手段と、前記計算した文字接続確率の値を蓄える確率テーブル記憶手段と、前記計算した文字接続確率を用いて文を単語単位に分割する文字列分割手段と、前記単語分割された文を出力する文書出力手段とを備えたことを特徴とする単語分割装置。

【請求項 10】 請求項 1 から 8 のいずれかに記載の単語分割方法または請求項 9 に記載の単語分割装置をソフトウェアにより実現したプログラムを記録した記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、電子計算機を利用した機械翻訳や大量文書検索、テキスト自動要約等を実施する自然言語処理システムの前処理・解析部における単語分割方法と装置に関し、特に、文を効率的に単語単位に分割できるようにした技術に関するものである。

【0002】

なお、以後の本発明の説明において、単語は文字の列（文字列）から構成されたもので、文字が組み合わさって意味を形成する単位とする。文は単語の列から

構成されているものであり、結果として文字列で表される。文書とは文が複数集まってまとまりを作った単位であるとする。

【0003】

【従来の技術】

日本語や中国語など単語を分けて書かない言語を膠着語という。膠着語では、言語の知識を持たない者がその字面だけ見ると、文は長い文字列であって、単語の境界をみつけることができない。機械翻訳や自動要約といった自然言語処理システムにおいては、その最初の段階として文の解析が必要になる。特に日本語のような膠着語では、単語への分割が最初の解析に相当する。

【0004】

また文書検索システムでは、例えば「今月の東京都議会」という文字列の「東京都議会」という語を単語の概念を使わずに文字列検索できるようにしてしまうと、「東京」で検索した場合でも「京都」で検索した場合でもヒットしてしまうことになるが、こうした検索ノイズを減らすためには、やはり文を単語に分割しておく必要がある。このような処理には、通常は形態素解析処理が行われる。形態素解析では、解析用の辞書を用いて文を単語へ分割処理するが、形態素解析の精度はこの辞書がどれだけ整っているかに依存する。

【0005】

一方、近年、文書中の文字列や文字の出現といったものを統計的に調べて、処理に必要な情報を得るという提案がなされている。これは例えば、既に単語に分割されている文書から、ある単語（または単語列）の次にどのような単語が出現しやすいかというものを確率として計算し、形態素解析の時にこの情報を使って解の候補を得るというものである。（参考文献：「単語と辞書」松本祐治 他 著 岩波書店 1997年）確率の計算には、単語Nグラムという単語N個組が使われる。NグラムはN-1個の単語の次にある単語が出現する確率を計算するもので、この確率計算はマルコフモデルとも呼ばれており、音声認識の単語推定などにも応用されている。ただし単語Nグラムは単語の接続可能性を計算するのであって、辞書にない単語を類推するものではない。

【0006】

この統計処理の考え方をを用いるものとして、単語のNグラムではなく、文字に着目した文字Nグラムがある。文字Nグラムは、 $(N-1)$ 個の文字列の後にどのような文字が続くかの確率を計算したものである。この文字Nグラムを応用し、文書中に出現する単語になりえそうな文字列の出現頻度を網羅的に調べて、その文字列前後の文字接続がどれほど散らばっているかを分散という尺度で計算することで単語や慣用句を収集する方法が特開平 9 - 1 3 8 8 0 1 号公報に開示されている。また、論文「正規化頻度による形態素解析の推定」(情報処理学会自然言語処理研究会 NL-113-3 1996)では、Nグラムの出現頻度を正規化計算することで、辞書を用いずに文を単語単位へ分割する手法が提案されている。

【0007】

辞書を用いない形態素解析は、特開平 1 0 - 3 2 6 2 7 5 号公報、特開平 1 0 - 2 5 4 8 7 4 号公報などにも開示されている。特開平 1 0 - 3 2 6 2 7 5 号公報は、文字Nグラムを使って、文字列の部分連鎖確率と単語分割点との関係をテーブルに記憶しておき、そのテーブルを使って単語分割を行うものである。テーブルの作成には、あらかじめ単語単位に分割された文(または文書)を用意し、その文書から計算機により自動学習を行わせる。特開平 1 0 - 2 5 4 8 7 4 号公報も同様に単語単位に分割された文書からあらかじめ学習をする必要がある。

【0008】

【発明が解決しようとする課題】

しかしながら、言語には常に新しい単語が生まれるものであるため、形態素解析用辞書は常にメンテナンスが必要である。また、対象とする文書によって単語の使われ方が違うこともあり、対象とする文書を変更する度に辞書を調整しなければいけない。そして、どれだけ注意していても、形態素解析において未知語、すなわち辞書に載っていない単語に遭遇する可能性は否定できず、未知語の出現により形態素解析の精度が低下することがある。

【0009】

辞書を使わない代わりに統計的処理を用いるものとして、前記の通り特開平 1 0 - 3 2 6 2 7 5 号公報記載の技術などがある。しかしながら、これらは事前に、単語単位に分割された文書を読ませることでシステムを訓練(自動学習)して

おく必要がある。単語単位に分割された文書を用意するためには、人手で文を分割しておくか、または既存の形態素解析システムの出力結果を用いる。だが人手で文を分割するのは多大なコストが必要であり、文書分野や時代ごとに大量の文を分割して用意することは難しい。そして時代の変化とともに変わっていく言語について、常に大量の分割済み文書を作成し続けなければならない、辞書の整備以上に大変な作業となる。また既存の形態素解析の出力結果を用いた場合、既存の形態素解析における解析失敗部分をそのまま学習してしまい、既存の形態素解析を越える精度は期待できなくなる。

【0010】

本発明は、上記の従来技術の課題を解決するためのものであり、基本的に辞書や単語単位に分割された大量の訓練用の文を必要とせず、文を単語へ分割することができる単語分割方法およびその方法を実施する装置を提供することを目的としている。

【0011】

【課題を解決するための手段】

上記目的を達成するために本発明は、単語に分割されていない文書から文字間の結合度を文字接続確率という形で計算し、この文字接続確率を使うことで文を単語単位へ分割するものである。これにより、辞書を使わず、また単語分割後文書を学習する必要もなく、文を単語に分割するという効果を奏するものである。

【0012】

【発明の実施の形態】

以下、本発明の実施の形態について説明する。まず前提となる言語の性質を説明する。文字の出現確率に注目してみる。一般に単語を構成する文字列は、全ての文字の組み合わせの単語が存在するわけではないので、文字の出現は等確率ではない。すなわちある言語の文字の種類をK種類とすると、もし単語を構成する文字が等確率で使われているなら、M文字からなる単語の種類はKのM乗個存在することになる。しかし、実際には語彙数はそれほど多くない。

【0013】

以後の説明では、膠着語として日本語を例にして説明する。日本語の日常生活

で通常使われる文字の種類は、約 6 千である。この数は、今日の一般的なコンピュータで扱える（JISで規定された）文字の種類数から類推したものである。いま日本語 2 文字の単語について考えると、もし全ての文字の組み合わせの単語が存在するなら、 6000 の 2 乗 = 3 千 6 百万の単語が存在することになる。日本語にはこの他に 3 文字や 4 文字の単語も存在するから、さらに多くの単語が存在することになる。しかし日本語の総語彙数は、たかだか数十万と考えられる。この数は、岩波書店の広辞苑等、日本語辞書の語彙数が 20 万から 30 万の間にあることから推測したものである。これについては、「自然言語処理」（長尾真編 岩波書店 1996 年出版）の第 2 章第 1 節「言語の統計」にも述べられている通りで、一般に文字の出現には偏りがあるものとされている。

【0014】

次に本発明の原理について説明する。ある文字「a」の後に別の文字「b」が続く確率は、もし前記の偏りがなければ、すなわちどんな文字も等確率で出現して単語を形成するなら、言語を構成する文字の種類 K の逆数（日本語の場合約 6 千分の 1）になる。しかし実際には偏りがあるのでそうはならない。具体例で説明すると、ある文字列「衆議院」が単語であったとすると、文字単位での接続確率を「衆」の後に「議」が続く条件付き確率 $P(\text{議} | \text{衆})$ とした時、この確率は日本語全体を調べてみるなら 6 千分の 1 より大きくなるはずである。同様に文字列「衆議」の後に文字「院」が続く条件付き確率 $P(\text{院} | \text{衆議})$ の場合は前 2 文字が与えられることから、さらに高い確率を示すはずである。一方で、存在しない単語（文字の組み合わせ）と思われる「衆び」などが出現する確率 $P(\text{び} | \text{衆})$ は、限りなく 0 に近くなるはずである。

【0015】

一方、文を構成する単語は、かなり自由な組み合わせが可能である。例えば「これは数学の本だ」「これは音楽の本だ」は両方とも文であるが、「数学」「音楽」の部分は自由な単語が接続できる。つまり単語を構成する文字列「これは」の後に別の単語の文字「数」が続く条件付き確率 $P(\text{数} | \text{これは})$ は、単語を構成する文字間の接続確率よりも低くなるはずである。この文字接続確率は、文字の間の結合度と解釈できる。そしてこれが計算できれば、それを基に文字列（文）

を単語単位へ分割することができる。

【0016】

文字間の接続確率は、文をある程度の量、つまりある程度の量の文書を集めることができれば、そこから統計的に調べて計算することができる。すなわち、文書データベースを構築するような状況ならば、データベースに登録する文書から文字接続確率を統計的に調べて計算することができる。この計算値は日本語全体について調べた場合の確率値とは違うものであろうが、近似できるものであり、しかもその確率値を調べた文書、あるいは類似の文書の分割に適用するのに適した性質を持つ。

【0017】

本発明では、以上の原理を用いることで、単語分割されていない文書から文字接続確率を調べ、その文字接続確率を使うことで、文書を辞書を用いることなく単語単位へと分割する。以下、本発明の実施の形態について図面を用いて説明する。

【0018】

(実施の形態1)

図1は本発明の実施の形態1における単語分割処理方法を示す流れ図である。図2は本実施の形態1における単語分割装置を示すブロック図であり、処理対象文書を電子化された状態で入力するための文書入力手段201と、入力した文書を蓄えておく文書データ蓄積手段202と、文書から文字接続確率を計算するための文字接続確率計算手段203と、計算した確率を記録しておくための確率テーブル記憶手段204と、計算した文字接続確率を使って文書を単語単位に分割するための文字列分割手段205と、処理結果の文書を出力する文書出力手段206とを備えている。

【0019】

以上のように構成された単語分割装置について、その動作を図1を用いて説明する。まず、ステップ101では、文書入力手段201から入力されたデータが、文書データ蓄積手段202に蓄えられる。ステップ102では、このデータから、文字間の接続確率を文字接続確率計算手段203が計算し、計算結果を確率

テーブル記憶手段 204 に蓄える。計算方法の詳細は後述する。ステップ 103 では、文字列分割手段 205 が、文書データ蓄積手段 202 に蓄えられたデータについて、確率テーブル記憶手段 204 に蓄えられた確率値を用いることで、文字間の接続確率を調べ、その確率が低い所で分割し、ステップ 104 で、分割された文を逐次、文書出力手段 206 から出力する。

【0020】

以上のように、本実施の形態 1 における単語分割装置は、処理対象文書から文字接続確率を計算し、計算した確率を使って対象文書を単語単位へ分割することができる。

【0021】

次に、図 1 のステップ 102 の詳細について説明する。本実施の形態 1 においては、文字 C_{i-1} と文字 C_i の間の文字接続確率は、文字列 $C_1 \dots C_{i-1}$ の後に次の文字 C_i が続く条件付き確率で表現することにする。これを式にすると次のように書ける。

【数 1】

$$P(C_i | C_1 C_2 \dots C_{i-1}) \quad (1)$$

【0022】

しかし、これは計算量が大きく記憶空間が大量に必要なことになる。(1) 式のような単語列や文字列の条件付き確率は、一般には N グラム ($N = 1, 2, 3, 4, \dots$) と良ばれる文字 N 個組で近似する。文字 N グラムによる条件付き確率とは、その $N - 1$ 個の文字列 $C_{i-N+1} \dots C_{i-1}$ という文字列が続いたという条件のもとで文字 C_i が出現する確率である。すなわち、 N グラムの 1 番目から $N - 1$ 番目の文字列が続いたという条件のもとで N 番目の文字が出現する確率である。これは次式 (2) のように書くことができる。

【0023】

【数 2】

$$P(C_i | C_{i-N+1} \dots C_{i-1}) \quad (2)$$

N グラムの確率は、文字列 $C_1 C_2 \dots C_m$ が調べようとするデータ中に出現する回数を $\text{Count}(C_1 C_2 \dots C_m)$ とすると、次式 (3) のように推定できる (参考文献: 「単語と辞書」 (松本祐治 他 著 岩波書店 1997 年))。

【数 3】

$$P(C_i | C_{i-N+1} \dots C_{i-1}) = \frac{\text{Count}(C_{i-N+1} \dots C_i)}{\text{Count}(C_{i-N+1} \dots C_{i-1})} \quad (3)$$

【0024】

なお、N グラム計算の時には、計算する文字列 (文) の前後に $N-1$ 文字の特殊な記号を付与して計算するのが一般的である。(参考文献: 同書) これは、一般の N グラムは文の先頭の文字の確率や、文の最後の文字の確率を、特殊記号を含めた N グラムにより計算するからである。N = 3 の例で説明するなら、「これは本だ」という文字列の 3 グラムを作成するためには、特殊記号をここでは # で表現することとして、「##これは本だ##」のような文字列を作成してから N グラムを作成する。すると「##こ」「#これ」「これは」「れは本」「は本だ」「本だ#」「だ##」の 7 つの 3 グラムを作ることになる。

【0025】

一方、本実施の形態 1 においては、一般的な N グラムの計算とは違い、計算する文字列の前後に $N-1$ 文字の特殊な記号を付与せず、計算する文字列の前にだけ $N-2$ 文字 (ただし $N-2 \geq 0$ とし、 $N=1$ の時は 0 とする) の特殊記号を付与する。これは文の後については、文の最後の文字の後には必ず単語として切れるのであるから、文の最後の文字とその後の特殊記号との接続確率を計算する必要がないからである。また文の前については、文の先頭が単語区切であることは自明であるので $N-1$ 個の特殊記号は必要なく、文頭 1 文字と次の 1 文字との接続確率を計算するために特殊記号 $N-2$ 個を含む N グラムを作成する必要があるのだ。

【0026】

先の「これは本だ」の 3 グラムの例で言うなら、文の先頭が単語区切であるの

は自明なので、「##こ」によって「##」と「こ」の接続確率を計算する必要は無いが、「#これ」によって「#こ」と「れ」の接続確率を計算する必要がある。文の前に必要な特殊記号の数は $N-2$ 個となる。同様に「本だ#」によって「本だ」と「#」の接続確率を計算する必要はないので、文末に必要な特殊記号の数は0個となる。

【0027】

式(3)を計算し、文字 N 個組とともに計算結果を図2の確率テーブル記憶手段204に記録することが図1のステップ102に相当する。確率テーブル記憶手段204は、例えば図4(d)のように、 N 文字組とその確率値が格納されるものであるが、文字組で検索しやすく記憶容量も小さくするために、適切な構造を用いて実現されているものとし、ここではその構造を限定しない。

【0028】

ステップ102の計算手順は、例えば図3に示す手順で実現できる。まず、ステップ301では、文書を構成する文ごとに、文の前に文頭を表現する特殊記号を $N-2$ 個付与する。ステップ302では、 $N-1$ グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 $N-1$ 個組について、それが何回出現しているかを調べた表を作成する。一般に N グラムの統計を調べる方法は、(参考文献:「言語情報処理」 長尾真 他 著 岩波書店 1998年)などに述べられているが、単純には文字の種類 K の N 乗を表現できるテーブルを用意し、そこに出現数をカウントして行くか、あるいは文書から全ての N 文字組みを取り出しそれをソートして同じものの出現回数をカウントすれば計算できる。

【0029】

ステップ303では、 N グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 N 個組について、それが何回出現しているかを調べた表を作成する。ステップ302と同様である。ステップ304では、 N グラム統計の夫々の文字 N 個組文字列について、その出現回数を X とし、同文字 N 個組文字列について、その1番目から $N-1$ 番目の文字列の出現回数をステップ301で作成した $N-1$ グラム統計から調べ、これを Y とし、 X/Y により式(3)の値を計算し、この値を確率テーブル記憶手段204に記録する。

【0030】

以上により式(3)の値が計算できるが、 $N-1$ グラムを作成しない方法も存在する。 N グラムは $N-1$ グラム文字列を含むことから、 N グラムを作っておけば $N-1$ グラムの出現頻度も簡単に計算できるからである。

【0031】

以下、文字接続確率の具体的な計算例を示す。全文書として文字列「a b a a b a」だけが与えられた場合を例とし、ここから文字接続確率を $N=3$ の N グラム(3グラム)で計算する。まずステップ301として、文(文字列)の前に文頭文末を表現する特殊記号を $N-2$ ($3-2=1$)個付与する。この様子を図4(a)に示す。特殊記号としてここでは#を付けているが、実際には文書に現れない記号を付けるものとする。次にステップ302として、2グラムの統計、すなわち文字2個組の出現回数を調べる。その結果が図4(b)のようになる。同様にステップ303として3グラムの、文字3個組の出現回数を調べ、図4(c)を得る。ステップ304として、図4(b)と図4(c)から、文字3個組についての式(3)の値を計算し、図4(d)を得る。以上が図1のステップ102の詳細説明である。

【0032】

次に、図1のステップ103の詳細について説明する。ステップ103はステップ102で計算した文字接続確率の表を使って、処理対象の文を構成する文字のそれぞれの部分の接続確率を調べ、分割をする過程である。その計算手順を図5に示す。本実施の形態1においては、 δ を閾値とし、この値はあらかじめ決められているものとする。まず、ステップ501で、文書から文の一つを選択する。次にステップ502で、図3のステップ301と同様に、文の前に文頭を表現する特殊記号を $N-1$ 個付与する。ステップ503では、ポインタを文の前に付けた特殊記号の一文字目に移動する。ステップ504では、ポインタ位置から始まる N 文字について、ステップ102で計算した文字接続確率を調べる。ステップ505では、もしその確率が、あらかじめ決められた閾値 δ 未満だったら、ポインタ位置を1文字目とした時の $N-1$ 文字目と N 文字目の間は単語分割点だったものと推定され、よってそこで分割を行う。もしその確率が閾値 δ 以上だったら

、そこは単語分割点ではないので分割を行わない。ステップ506で、ポインタを一文字進める。ステップ507で、ポインタを1文字目とした時のN文字目が文末文字を越える場合、文は終了したものとして、ステップ508へ、そうでなければステップ504へジャンプする。ステップ508で、文書から次の文を選択する。ステップ509で、次の文が無ければ終了、そうでなければステップ502へ進む。

【0033】

以上により分割点を発見する。以下、具体的計算例を、先に示した図4の文字列「a b a a b a」のN=3の場合で示す。既に図4(d)の文字3個組の接続確率は計算されているものとする。閾値として $\delta = 0.7$ が与えられているものとする。まず、この例の場合では文が一つしかないので、ステップ501で「a b a a b a」が選択され、ステップ502で文の前に特殊文字が付けられることで図4(a)と同じ状態になる。次にステップ503でポインタを移動させた状態が図4(e)である。ここから3文字、すなわち「# a b」の確率を図4(d)のテーブルで探すと1.0であり、これは閾値 $\delta = 0.7$ より大きいので、「# a」と「b」の間は分割されない。以下同様にステップ504からステップ506を繰り返すことで、文字列のそれぞれの点での接続確率が調べられ、この値をもって単語分割点を決定することができる。この様子を図4(f)に示す。この例では、単語分割された結果は「a b a / a b a」となる。

【0034】

もう一つ別の例として、同様に日本語の単純な文字列「にわにわにわ」（庭には二羽）を計算したのが図6である。文頭特殊記号を付与したものが図6(a)であり、2グラムと3グラムの出現回数はそれぞれ図6(b)、(c)に計算され、そこから3グラムの文字接続確率は図6(d)のようになる。これを元の文字列にあてはめていくと、図6(e)のようになり、閾値 $\delta = 0.7$ とすることで、結果として「にわ / にわ / にわ」と分割される。

【0035】

上記2例とも文中の文字の種類が少ない例を示した。これは計算例として示すために、非常に短い文から確率計算をしたためである。日本語のように文字種が

多い場合は、さらに多くの学習用（確率計算用）の文が必要である。

【0036】

新聞データ約1千万文字からなる文書（文の集合）で文字間の接続確率を計算した例について、その一部を示したものが図7である。この計算結果を使い、文字列「利用者の減少と反比例するように」を計算した結果が図8である。図8では閾値 $\delta = 0.07$ で分割点を決定している。

【0037】

本実施の形態1においては、分割処理対象文書自身から文字接続確率を計算し、その確率を使って同じ分割処理対象文書を分割した。この方法は対象文書に出現する文字の組み合わせの確率全てを計算できるという点で合理的である。

【0038】

なお、本発明は、処理対象文書自身だけから文字接続確率を計算するというものに限定されるものではない。まとまった文書から文字接続確率を計算しておき、それを使って別の文書を分割することも可能である。これは漸増的な文書データベースにおいて有効である。この場合は分割対象文書に出現する文字の組み合わせが、確率を計算（学習）した文書に出現していない可能性も否定できないが、これらはNグラム平滑化の問題として（参考文献：「単語と辞書」松本祐治 他 著 岩波書店 1997年）などに記述されている方法で対応できる。

【0039】

以上のように、本実施の形態1では、ステップ101で入力された文書からステップ102で文字間の接続確率を計算し、この接続確率を使ってステップ103で該文書のそれぞれの文字の接続確率を調べることで単語分割を行い、ステップ104で結果を出力することで、辞書を使わない単語分割を行うことが可能になり、その実用的効果は大きい。

【0040】

（実施の形態2）

次に、本発明の第2の実施の形態について説明する。本実施の形態2における単語分割装置の構成は、図2に示した実施の形態1と同じものである。また動作の概要は、実施の形態1と略同様であるが、計算方法として別のものを用い、

図 1 のステップ 1 0 2、およびステップ 1 0 3 の手順が変更されているので、その詳細について説明する。

【0 0 4 1】

上記した実施の形態 1 の説明では、文字列接続確率の計算には N グラムを用いることで、文字列 $C_{i-N+1} \dots C_{i-1}$ が出現したという条件のもとで文字 C_i が出現する確率を使った（式（2）参照）。例えば文中に出現する「a b c d e f」の「a b c」と「d e f」の間の接続確率を計算するために、「a b c」という文字列が出現した場合に次が「d」である確率を使ったのである。これは既存の技術である N グラム方法を転用して用いたからである。N グラムはもともと、単語の接続の確からしさや文字の接続の確からしさを計算し、文全体として正しいかを判断するためのものである。または、いままでに出現した単語列や文字列から次の単語や文字を予想するものである。従って、本来は次式（4）という確率式で計算するものを、式（5）で近似した場合の product 記号 Π の中の項であった。つまり全ての項を掛け合わをせる形で使うことが前提だったので、文字列 $C_{i-N+1} \dots C_{i-1}$ が出現したという条件のもとで文字 C_i が出現する確率というように、条件の部分が複数文字（文字列）で、その後特定の 1 つの文字が出現する確率で扱えたのである。

【0 0 4 2】

【数 4】

$$\prod_{i=1}^m P(w_i | w_1 w_2 \dots w_{i-1}) \quad (4)$$

【数 5】

$$\prod_{i=1}^m P(w_i | w_{i-N+1} \dots w_{i-1}) \quad (5)$$

【0 0 4 3】

しかし、本発明は文字接続確率を、単語内での文字接続か、単語間の文字接続かを判別するために用いる。従って、本実施の形態 2 では、文字 C_{i-1} と文字 C_i の接続確率を、ある文字列が出現したという条件におけるある 1 文字の出現確率で表現するのではなく、ある文字列が出現したという条件におけるある文字列の出現確率を計算するという方法にする。

【0044】

形式的に表現するなら、長さ n 個の文字列 $C_{i-n} \dots C_{i-1}$ が出現したという条件のもとで長さ 1 個の文字列 C_i が出現する確率を計算するよりも、長さ n 個の文字列 $C_{i-n} \dots C_{i-1}$ が出現したという条件のもとで長さ m 個の文字列 $C_i \dots C_{i+m-1}$ が出現する確率を計算する。この確率を式 (2) と同様に書くならば、次式 (6) のようになる。

【0045】

【数 6】

$$P(\underbrace{C_i \dots C_{i+m-1}}_{m \text{ 個}} | \underbrace{C_{i-n} \dots C_{i-1}}_{n \text{ 個}}) \quad (6)$$

【0046】

例えば文中に出現する「a b c d e f」の「a b c」と「d e f」の間の接続確率を計算するために、「a b c」という文字列が出現した場合に、次が「d e f」である確率を使うのである。これは $n = 3$ 、 $m = 3$ の例である。式 (6) は $m = 1$ の場合が上記した実施の形態 1 に相当する。

【0047】

また上記した実施の形態 1 が、文の先頭側にある文字列から次の文字列への接続確率を求める、つまり前から後へ進む順方向の確率の計算と考える時、式 (6) の $n = 1$ 、 $m > 1$ という条件は、後から前への接続確率の計算に近似でき、逆方向の確率の計算に相当する。例えば文中に出現する「a b c d e f」の「a b c」と「d e f」の間の接続確率を計算するために、 $n = 1$ 、 $m = 3$ なら、文字「c」の後に文字列「d e f」が出現する確率ということになるが、文字列「d

e f」が出現した場合にその前が「c」である確率に近似できる。これは逆方向の文字接続確率の計算に相当する。ところが式（6）の計算のためには、 $n + m$ グラムの統計を取る必要がある。しかし、 $n \geq 2$ かつ $m \geq 2$ とすると、4 グラム（またはそれ以上）の文字組の統計計算が必要になり、非常に大きな記憶空間を必要とする。

【0048】

そこで本実施の形態2では、式（6）の計算を次式（7）で近似する方法を提案する。

【数7】

$$P(C_i | C_{i-n} \dots C_{i-1}) \times P(C_{i-1} | C_i \dots C_{i+m-1}) \quad (7)$$

【0049】

式（7）は、 n 個の文字列が出現した後に特定の文字が出現する順方向の確率である第1項と、 m 個の文字列が出現する前に特定の文字が出現する逆方向の確率である第2項の積である。項と文字列の関係を図14に示す。例えば文中に出現する「a b c d e f」の「a b c」と「d e f」の間の接続確率を計算するために、「a b c」の後に「d」が出現する確率（順方向の第1項）と、「d e f」が出現した時にその前が「c」である確率（逆方向の第2項）の積を取ることを意味する。

【0050】

式（7）の確率値は、第1項は $n + 1$ グラム、第2項は $m + 1$ グラムの統計を取ればよく、次式（8）により計算（推定）できる。

【数 8】

$$\underbrace{\frac{\text{Count}(C_{i-n} \dots C_i)}{\text{Count}(C_{i-n} \dots C_{i-1})}}_{\text{第 1 項}} \times \underbrace{\frac{\text{Count}(C_{i-1} \dots C_{i+m-1})}{\text{Count}(C_i \dots C_{i+m-1})}}_{\text{第 2 項}} \quad (8)$$

【0 0 5 1】

式 (8) を計算し、文字 $n + 1$ 個組、および文字 $m + 1$ 個組とともに計算結果を、図 2 の確率テーブル記憶手段 2 0 4 に記録することが、本実施の形態 2 における図 1 のステップ 1 0 2 に相当する。よって、図 2 の確率テーブル記憶手段 2 0 4 は、文字 $n + 1$ 個組用と文字 $m + 1$ 個組用の 2 つのテーブルを持つことになる。 $n \neq m$ の場合、この計算手順は、例えば図 9 に示す手順で実現できる。

【0 0 5 2】

まずステップ 9 0 1 で、文書を構成する文ごとに前後に文頭文末を表現する特殊記号を、文頭には $n - 2$ 個、文末には $m - 2$ 個付与する。本実施の形態 2 では、後から前への接続確率も計算するため、文の後にも特殊記号を付与する必要がある。次にステップ 9 0 2 で、 n グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 n 個組について、それが何回出現しているかを調べた表を作成する。ステップ 9 0 3 では、 $n + 1$ グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 $n + 1$ 個組について、それが何回出現しているかを調べた表を作成する。ステップ 9 0 4 では、 $n + 1$ グラム統計の夫々の文字 $n + 1$ 個組文字列について、その出現回数を X とし、同文字 $n + 1$ 個組文字列について、その 1 番目から n 番目の文字列の出現回数をステップ 9 0 2 で作成した n グラム統計から調べ、これを Y とし、 X / Y により式 (8) の第 1 項の値を計算し、計算結果を確率テーブル記憶手段 2 0 4 の文字 $n + 1$ 個組 (第 1 項の確率値) の部分に記録する。次に、ステップ 9 0 5 で、 m グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 m 個組について、それが何回出現しているかを調べた表を作成する。ステップ 9 0 6 では、 $m + 1$ グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 $m + 1$ 個組について、それ

が何回出現しているかを調べた表を作成する。ステップ 907 では、 $m+1$ グラム統計の夫々の文字 $m+1$ 個組文字列について、その出現回数を X とし、同文字 $m+1$ 個組文字列について、その 2 番目から $m+1$ 番目の文字列の出現回数をステップ 905 で作成した m グラム統計から調べ、これを Y とし、 X/Y により式 (8) の第 2 項の値を計算し、計算結果を確率テーブル記憶手段 204 の文字 m 個組 (第 2 項の確率値) の部分に記録する。

【0053】

以上が $n \neq m$ の場合である。 $n = m$ の場合、図 2 の確率テーブル記憶手段 204 は、文字 n 個組用だけを用意すればよく、その構造は図 12 (d) のように n 文字組とその第 1 項の確率、第 2 項の確率を記録するものとなる。そして $n = m$ の場合、計算手順は例えば図 10 に示すように図 9 よりも簡略化できる。

【0054】

まずステップ 1001 で、文書を構成する文ごとに前後に文頭文末を表現する特殊記号を $n-2$ 個ずつ付与する。次にステップ 1002 で、 n グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 n 個組について、それが何回出現しているかを調べた表を作成する。ステップ 1003 では、 $n+1$ グラム統計を作成する。すなわち、対象文書の中に出現した全ての文字 $n+1$ 個組について、それが何回出現しているかを調べた表を作成する。ステップ 1004 では、 $n+1$ グラム統計の夫々の文字 $n+1$ 個組文字列について、その出現回数を X とし、同文字 N 個組文字列について、その 1 番目から n 番目の文字列の出現回数をステップ 1002 で作成した n グラム統計から調べ、これを Y とし、 X/Y により式 (8) の第 1 項の値を計算し、計算結果を確率テーブル記憶手段 204 の確率値第 1 項の部分に記録する。また、ステップ 1005 では、 N グラム統計の夫々の文字 $n+1$ 個組文字列について、その出現回数を X とし、同文字 N 個組文字列について、その 2 番目から $n+1$ 番目の文字列の出現回数をステップ 1002 で作成した n グラム統計から調べ、これを Y とし、 X/Y により式 (8) の第 2 項の値を計算し、計算結果を確率テーブル記憶手段 204 の確率値第 2 項の部分に記録する。

【0055】

以上により式(8)の値が計算できる下地が整ったが、まだ式(8)の値そのものは求めてなく、実際の値は次の分割過程で計算する。以下では図1のステップ103に対応する処理について詳細を説明する。図1のステップ103はステップ102で計算した文字接続確率の表を使って、処理対象の文を構成する文字のそれぞれの部分の接続確率を計算し、分割をする過程である。その計算手順を、 $n \neq m$ の場合について図11に示す。

【0056】

まずステップ1101で、文書から文を一つ選択する。次にステップ1102で、図10のステップ1001と同様に、文の前後に文頭文末を表現する特殊記号を、文頭には $n-2$ 個、文末には $m-2$ 個付与する。ステップ1103では、ポインタを文の前に付けた特殊記号の一文字目に移動する。ステップ1104では、ポインタ位置から始まる $n+1$ 文字について、確率テーブル記憶手段204の確率値第1項から検索し、それをポインタ位置を1文字目とした時の n 文字目と $n+1$ 文字目の間の確率値第1項として記録する。ただし文字と文頭文末特殊記号の間の接続確率は0とする。ステップ1105では、同様にポインタ位置から始まる $m+1$ 文字について、確率テーブル記憶手段204の確率値第2項から検索し、それをポインタ位置を1文字目とした時の1文字目と2文字目の間の確率値第2項として記録する。ただし文字と文頭文末特殊記号の間の接続確率は0とする。ステップ1106で、ポインタを一文字進める。ステップ1107で、ポインタが文末の文字を指していたら、文は終了したものとして、ステップ1108へ、そうでなければステップ1104へジャンプする。ステップ1108では、各文字間について、各文字間に記録された確率値第1項と確率値第2項の積を取り、式(8)の値を計算する。それがあらかじめ決められた閾値 δ 未満だったら、そこで分割を行う。もしその確率が閾値 δ 以上だったら、そこは単語分割点ではないので分割を行わない。ステップ1109で、文書から次の文を選択する。ステップ1110では、次の文が無ければ終了し、そうでなければステップ1102へ進む。以上により分割点を発見する。 $n = m$ の場合も、同様である。

【0057】

以下、具体的計算例を示す。全文書として文字列「仕事は仕事」だけが与えら

れた場合を例とし、ここから文字接続確率を $n=m=2$ の $(n+1)$ グラム (3 グラム) で計算する。まずステップ 1001 として、文 (文字列) の前後に文頭文末を表現する特殊記号を $n-2$ ($3-2=1$) 個ずつ付与する。この様子を図 12 (a) に示す。特殊記号としてここでは # を付けているが、実際には文書に現れない記号を付けるものとする。次にステップ 1002 として、2 グラムの統計、すなわち文字 2 個組の出現回数を調べる。その結果が図 12 (b) のようになる。同様にステップ 1003 として 3 グラムの、文字 3 個組の出現回数を調べ、図 12 (c) を得る。またステップ 1004 として、図 12 (b) と図 12 (c) から、文字 3 個組についての式 (8) の第 1 項の値を計算し、図 12 (d) の第 1 項の部分を得る。またステップ 1005 として、図 12 (b) と図 12 (c) から、文字 3 個組についての式 (8) の第 2 項の値を計算し、図 12 (d) の第 2 項の部分を得る。

【0058】

注意しなければならないのは、図 12 (d) は同じ箇所の文字接続確率を記録したものではない。図 12 (d) の表の第 2 行目「仕事は」について、第 1 項の部分に記録された確率は「仕事」と「は」の接続確率第 1 項であり、第 2 項の部分に記録された確率は「仕」と「事は」の接続確率第 2 項である。以上により確率値のテーブルができたので、次に図 11 の処理に進む。

【0059】

ステップ 1101 として「仕事は仕事」が選択され、ステップ 1102 で前後に特殊文字が付けられることで図 12 (a) と同じ状態になる。ステップ 1103 からステップ 1105 で図 12 (e) を得る。第 2 項は文頭特殊記号との間なので 0 となる。同様にステップ 1103 から 1106 までの繰り返しで図 12 (f) を得る。ステップ 1108 により各文字間の接続確率が計算され、図 12 (f) の接続確率の部分を得る。あらかじめ決められた閾値 δ (ここでは閾値 $\delta = 0.6$ とする) 未満の部分で分割を行うと、図 12 (f) に示す通り、「仕事／は／仕事」という分割結果を得る。

【0060】

以上のように、本実施の形態 2 では、文字 n 個の後に文字 m 個が続く確率を近

似式（８）で計算することで、より正確に、辞書を使わない単語分割を行うことが可能になり、その実用的効果は大きい。

【 0 0 6 1 】

なお、本実施の形態 1 および 2 では、閾値 δ はあらかじめ決められたものとして扱ってきたが、確率の値を計算した後、単語分割結果が望むべく平均単語長を満すように動的に決めてもよい。すなわち、図 1 3 に示すように、閾値 δ が大きければ平均単語長は長くなり、閾値 δ を小さくすれば平均単語長は短くなる。分割結果で調整しながら δ を文書ごとに決めるようにすれば、適切な値が取れるようになる。

【 0 0 6 2 】

また、閾値 δ は一率として扱ってきたが、何らかの基準により複数設定してもよい。日本語の場合は本来、平仮名部分の平均単語長は漢字部分の平均単語長よりも短い傾向にある。これは平仮名が助詞などの一文字の単語を含むからである。その一方で片仮名部分は外国語の発音を表記したものが多いことから平均単語長は長い。よって閾値 δ を文字種（漢字・平仮名・片仮名）により複数設定してもよい。

【 0 0 6 3 】

また、日本語の場合、単語分割点は文字種（漢字・平仮名・片仮名）の変化点にあることが多い。よって文字種の変化点の閾値 δ を他の部分より下げるなどしてより適切な値に調整してもよい。

【 0 0 6 4 】

また、本発明の実施の形態 1 および 2 では、文頭や文末は必ず単語分割点であるとして説明してきたが、この他に句読点の前後、カッコや記号の前後も単語の分割点とみなしてよく、その部分の確率計算を省略することが可能である。あるいは N グラム統計作成において、句読点や記号の前後も文の切れ目として計算してよい。すなわち、本発明第 2 の実施例で用いた「仕事は仕事」の 3 グラムは、文頭文末特殊記号を付与することで「#仕事は仕事#」という文字列を作って計算した。これが「仕事は、仕事」だった場合、文は 2 つであるとみなし、「#仕事は#」と「#仕事#」の 2 つの文について計算するようにしてもよい。

【0065】

また、本発明の実施の形態2では、式(6)の計算を式(7)のように第1項と第2項の積の形で近似したが、積以外の、例えば平均を取るなどの計算方法を用いてもよい。

【0066】

【発明の効果】

以上のように、本発明では、処理対象文書の中から文字接続確率を計算し、その確率を処理対象文書に当てはめることで単語分割できる場所を発見して分割するものであり、これにより、辞書を使わずにテキストを単語に分割するという効果を奏するものである。従って、本発明は単語の分割のために辞書を用意する必要がないことから、日々生まれ続ける新しい語や語法のために、辞書を整備したり各種パラメータを整備する必要もない。また、辞書を持たないことから、本発明の方法により実現されたプログラムを格納した記録媒体は非常に小さくて済む。同時にパーソナルコンピュータなどの処理能力に限界のある環境下においても機能する。従って、その実用的効果は非常に大きい。

【図面の簡単な説明】

【図1】

本発明の第1の実施の形態における単語分割方法の動作を示すフロー図

【図2】

本発明の第1および第2の実施の形態における単語分割装置の構成を示すブロック図

【図3】

本発明の第1の実施の形態における文字接続確率の計算手順を示すフロー図

【図4】

本発明の第1の実施の形態における単語分割計算例を示す概念図

【図5】

本発明の第1の実施の形態における分割過程の計算手順を示すフロー図

【図6】

本発明の第1の実施の形態における単語分割計算例を示す概念図

【図 7】

本発明の第 1 の実施の形態における単語分割方法で新聞記事データの文字接続確率を計算した例の一部を示す数値図

【図 8】

本発明の第 1 の実施の形態における単語分割方法で新聞記事データの分割をした例を示す概念図

【図 9】

本発明の第 2 の実施の形態における文字接続確率の計算手順について、 n と m が違う場合の計算手順を示すフロー図

【図 1 0】

本発明の第 2 の実施の形態における文字接続確率の計算手順について、 n と m が同じ場合の計算手順を示すフロー図

【図 1 1】

本発明の第 2 の実施の形態における分割過程の計算手順を示すフロー図

【図 1 2】

本発明の第 2 の実施の形態における単語分割計算例を示す概念図

【図 1 3】

本発明における閾値と平均単語長の関係を示す概念図

【図 1 4】

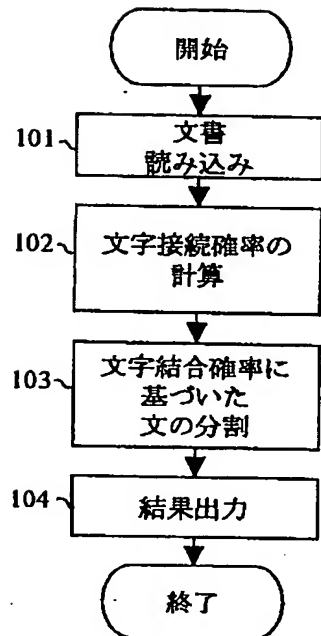
本発明の第 2 の実施の形態における式 (7) の関係を示す模式図

【符号の説明】

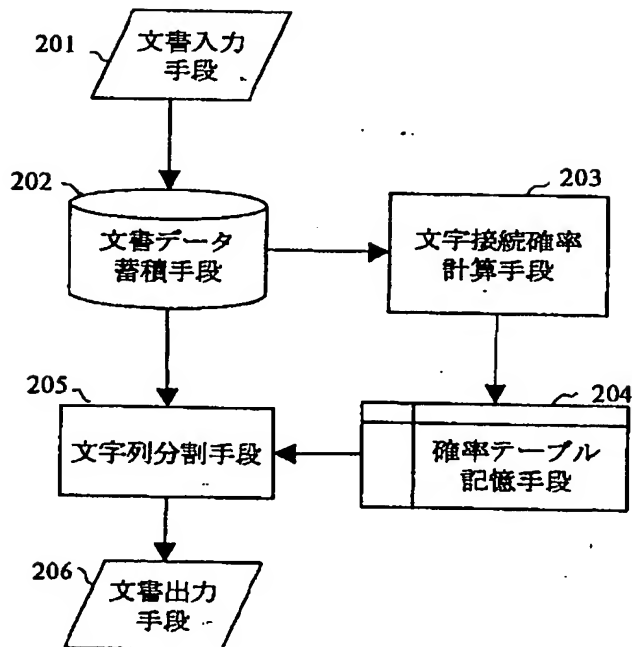
- 2 0 1 文書入力手段
- 2 0 2 文書データ蓄積手段
- 2 0 3 文字接続確率計算手段
- 2 0 4 確率テーブル記憶手段
- 2 0 5 文字列分割手段
- 2 0 6 文書出力手段

【書類名】 図面

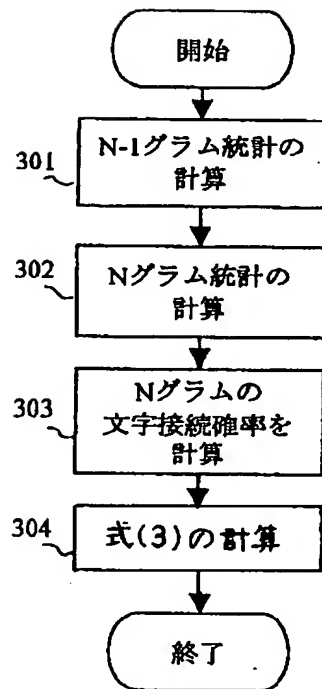
【図 1】



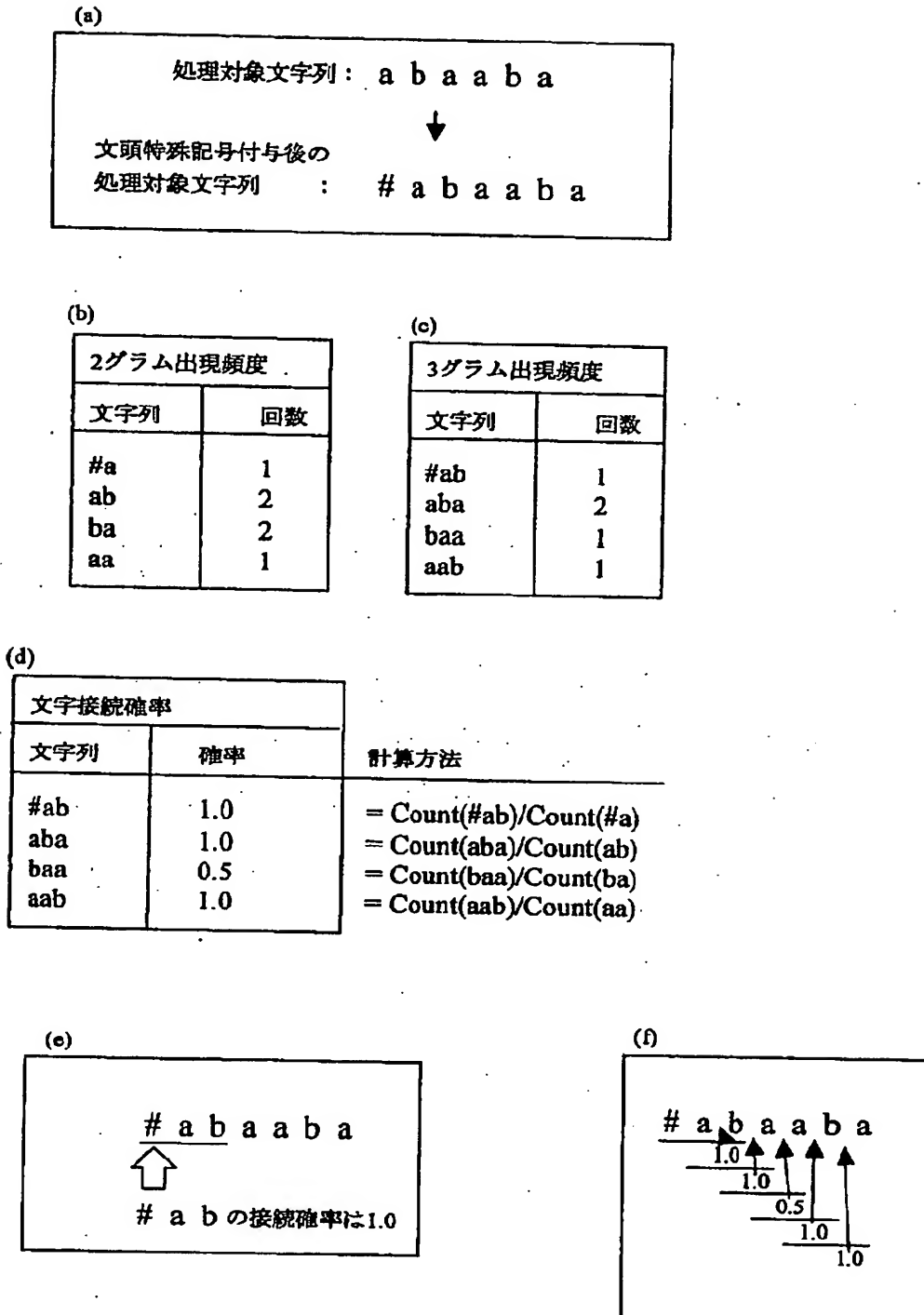
【図 2】



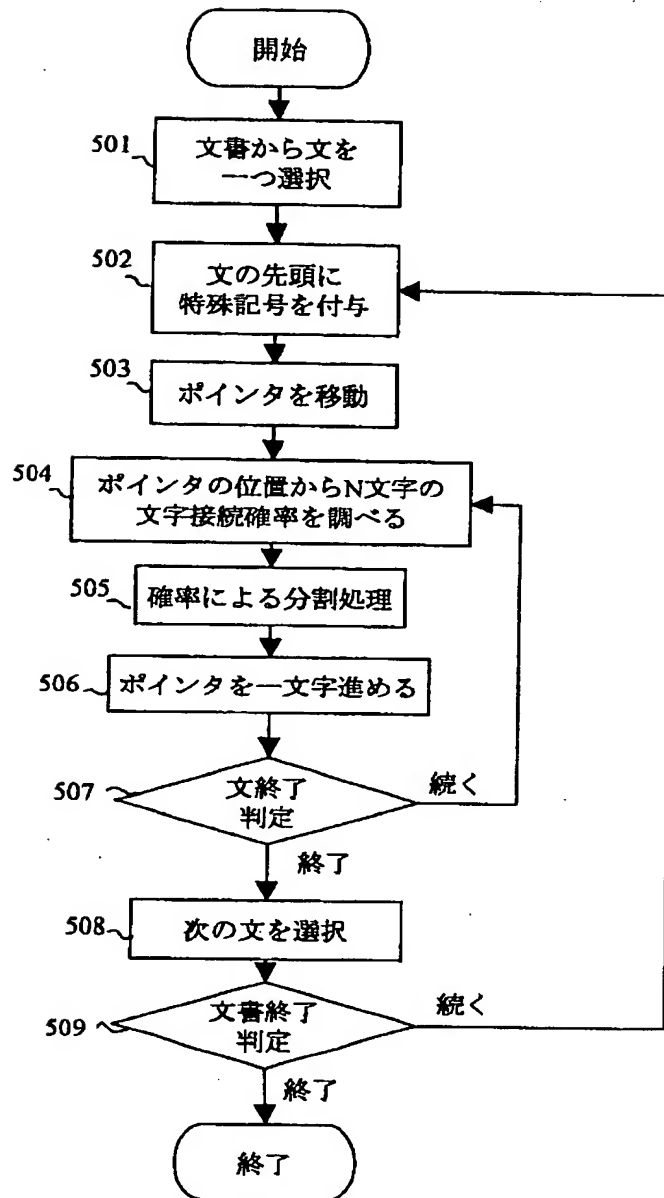
【図 3】



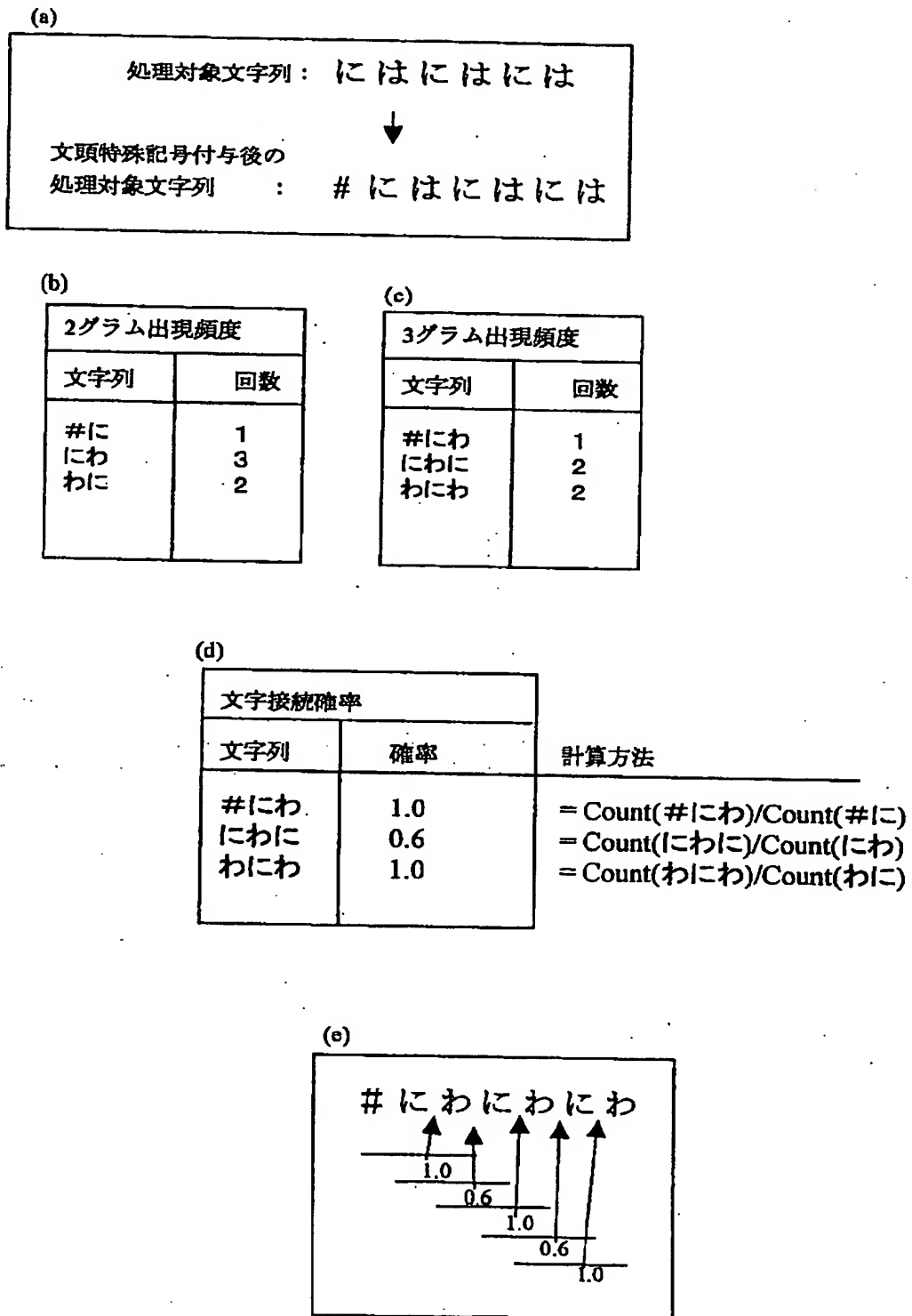
【図 4】



【図 5】



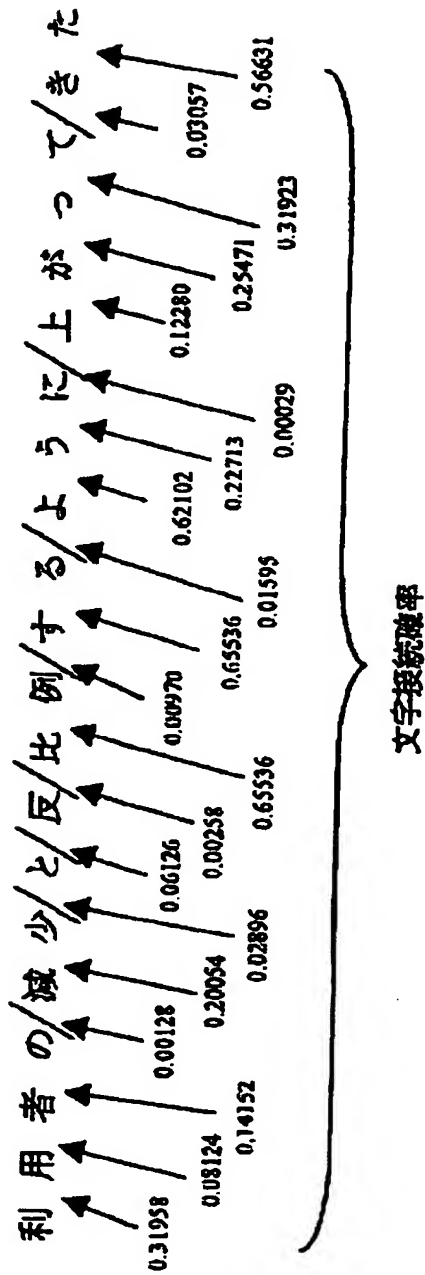
【図 6】



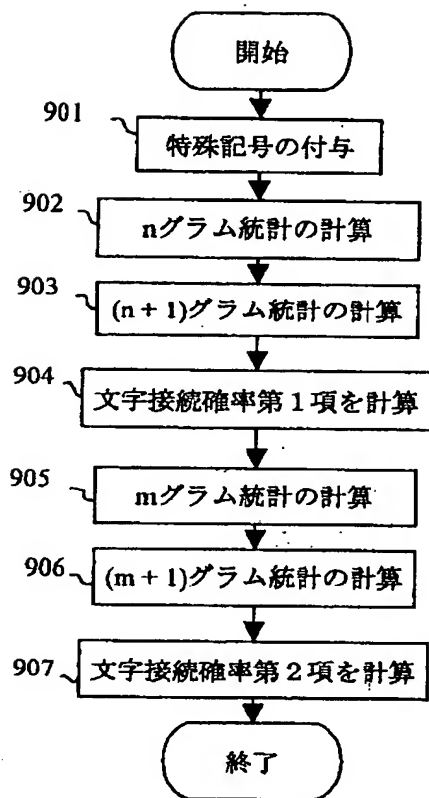
【図 7】

文字接続確率	
文字列	確率
ようち	0.00016
ようっ	0.00011
ようつ	0.00005
ようて	0.00011
ようで	0.01938
ようと	0.08398
ような	0.13665
ように	0.22713
よくあ	0.01593
よくい	0.00612
よくえ	0.00024
よくお	0.00269
よくか	0.00171

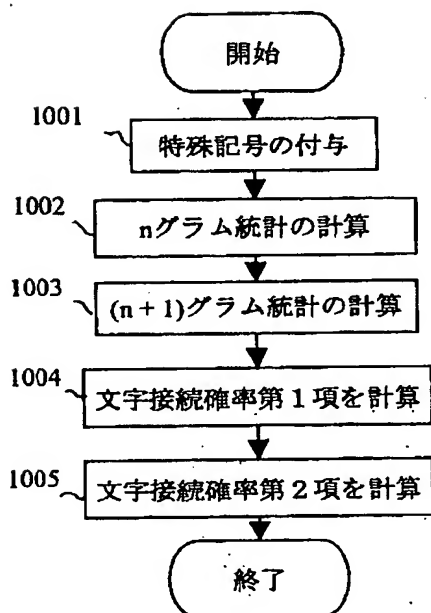
【図 8】



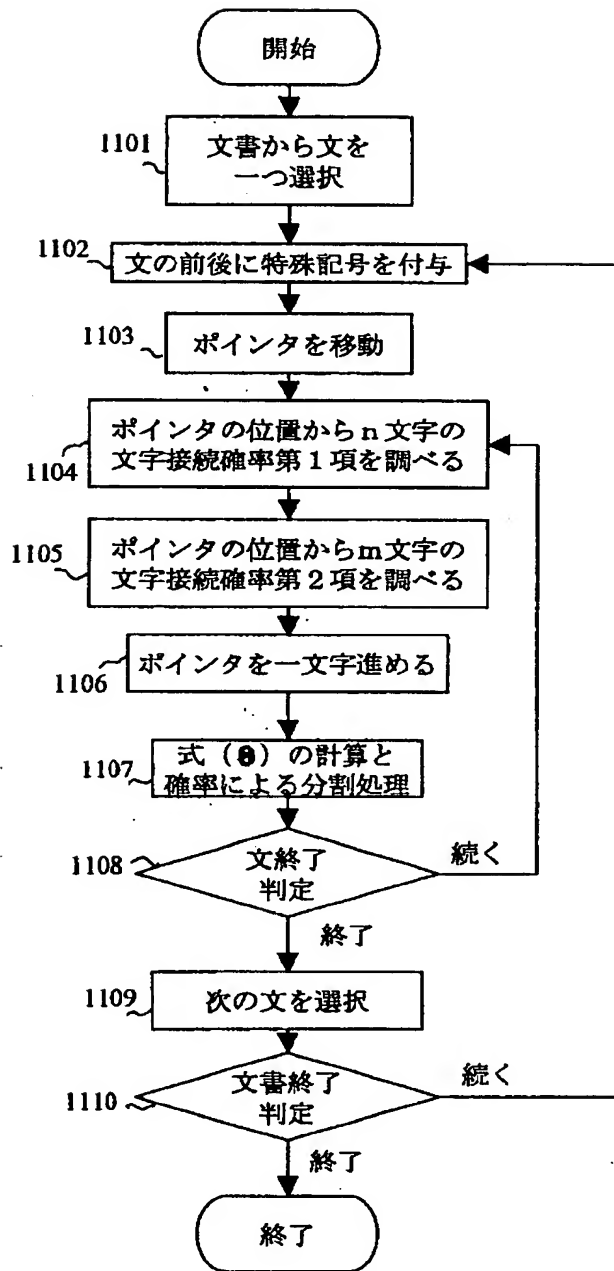
【図 9】



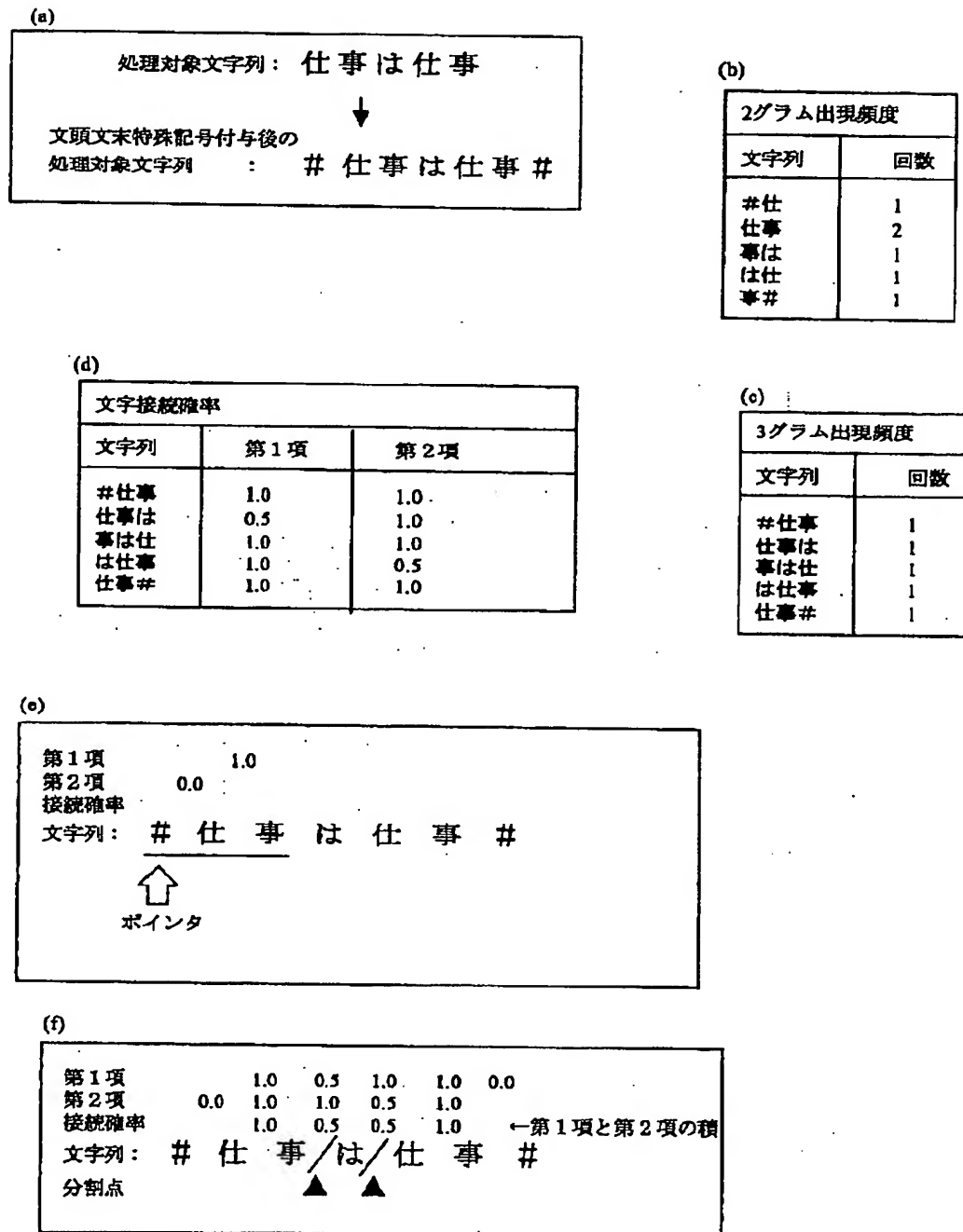
【図 1 0】



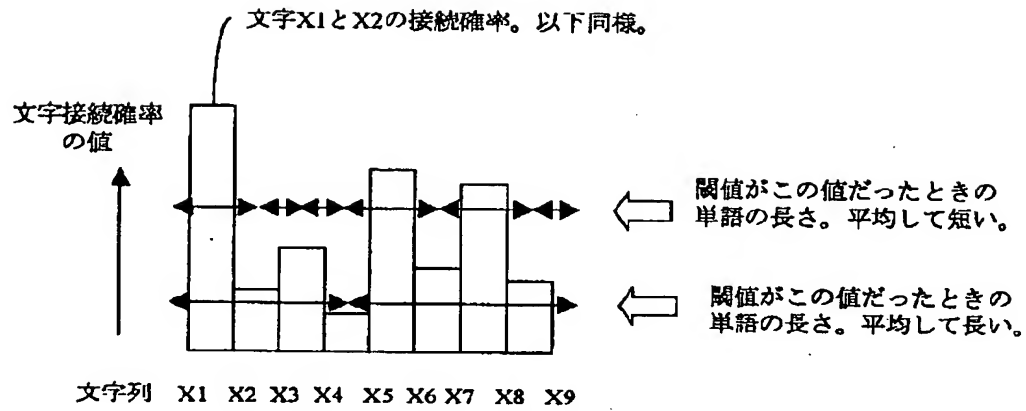
【図 1 1】



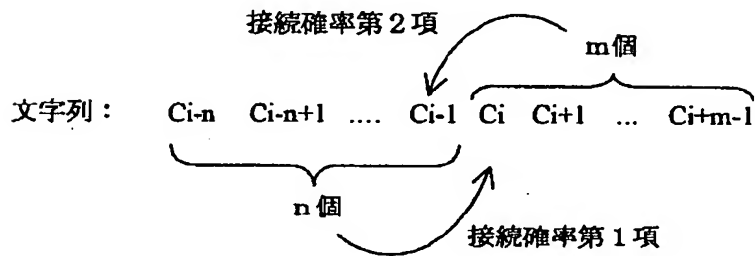
【図 12】



【図 1 3】



【図 1 4】



【書類名】 要約書

【要約】

【課題】 電子計算機を利用した自然言語処理システムにおいて、辞書や学習用単語分割済文を必要としない単語分割を実現することを目的とする。

【解決手段】 入力された単語分割されていない文書から、文字結合度としての文字接続確率を計算し、テーブルに記録する。この文字接続確率を用いて入力された文書を調べ、接続確率の低い部分で文書を分割し、出力する。

【選択図】 図 1

出 願 人 履 歷 情 報

識別番号 [000005821]

1. 変更年月日 1990年 8月28日

[変更理由] 新規登録

住 所 大阪府門真市大字門真1006番地
氏 名 松下電器産業株式会社